

Research at the Scalable Computing Software Laboratory

Xian-He Sun

Zhiling Lan

Department of Computer Science
Illinois Institute of Technology

www.cs.iit.edu/~scs/

Fermi Lab Presentation

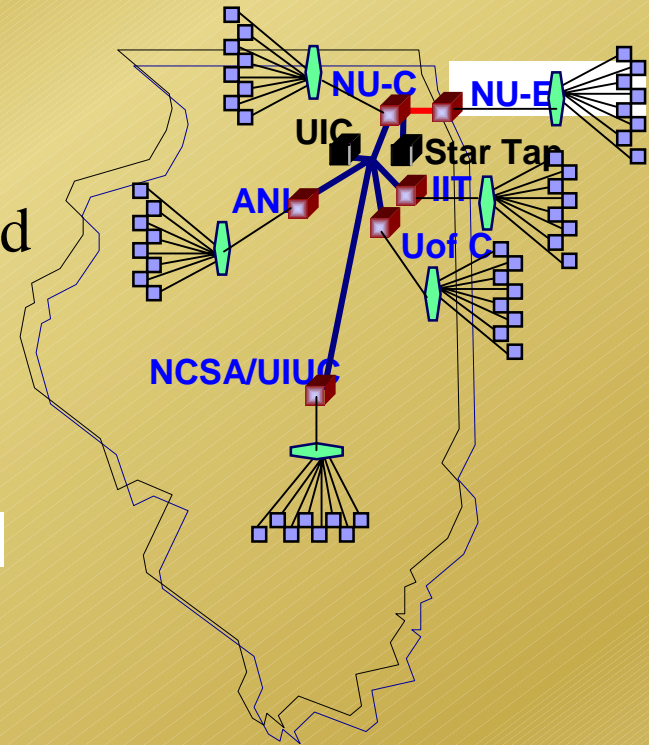
SCS Computing Infrastructure



Parallel Computers

Distributed
Optical Testbed

I-WIRE
OMNI



Pervasive Computing
Environments

X. Sun, Feb. 2004

SCS

Outline

- High Performance Computing
 - Scalable numerical kernel solvers
 - Performance optimization
- Distributed Computing
 - Mobility and mobility of legacy code
 - Performance prediction and task scheduling
- Pervasive Computing (not cover)
- Application
- Conclusion

Scalable Numerical Algorithms

- Motivation
 - Parallel codes have been developed during last decade
 - The performances of many codes suffer in a scalable computing environment
- Achievement
 - Scalable tridiagonal solvers
 - Fast and high-order Poisson solvers
 - Iterative Helmholtz equation algorithms
 - Domain decomposition methods

Sample: The PDD algorithm & its applications

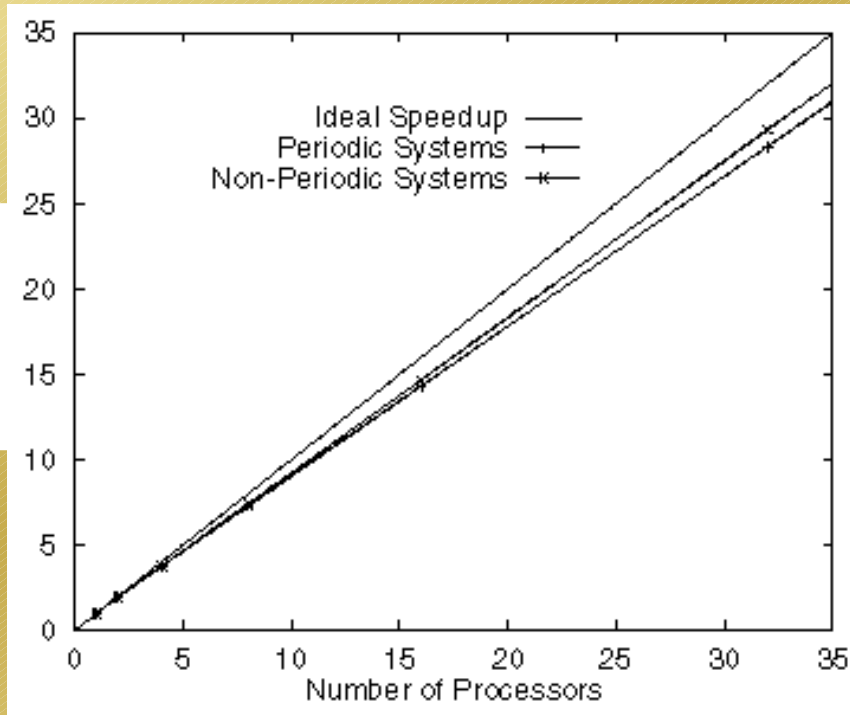
The Parallel Diagonal Dominant (PDD)

The Reduced System

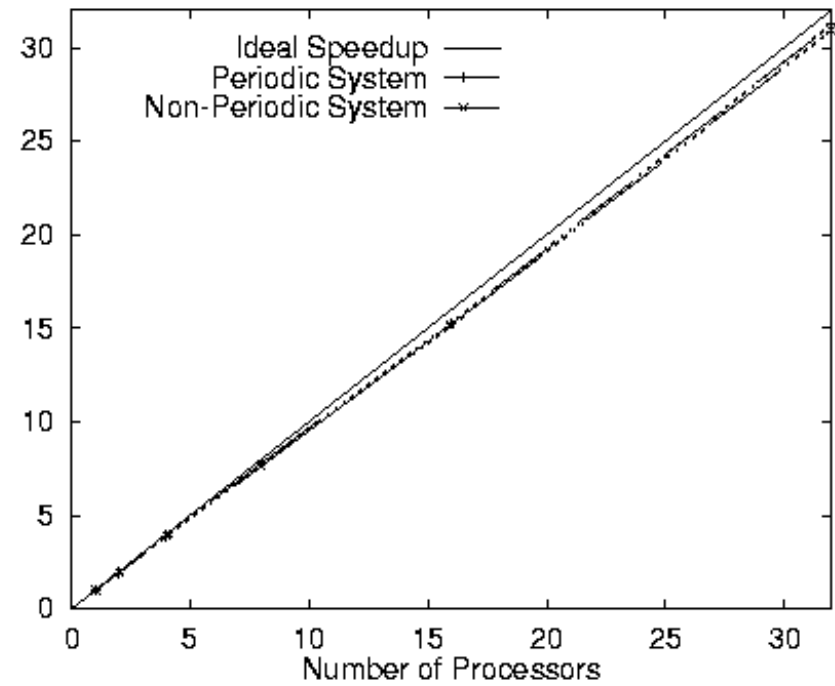
$$\begin{bmatrix} \mathbf{I} & \text{red} & & \\ \text{red} & \mathbf{I} & \text{red} & \\ & & \text{red} & \text{red} \\ & & & \mathbf{I} \end{bmatrix} \begin{pmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{pmatrix} = \begin{pmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{pmatrix}$$

Generally needs global communication, Decay for diagonal dominant systems

speedup



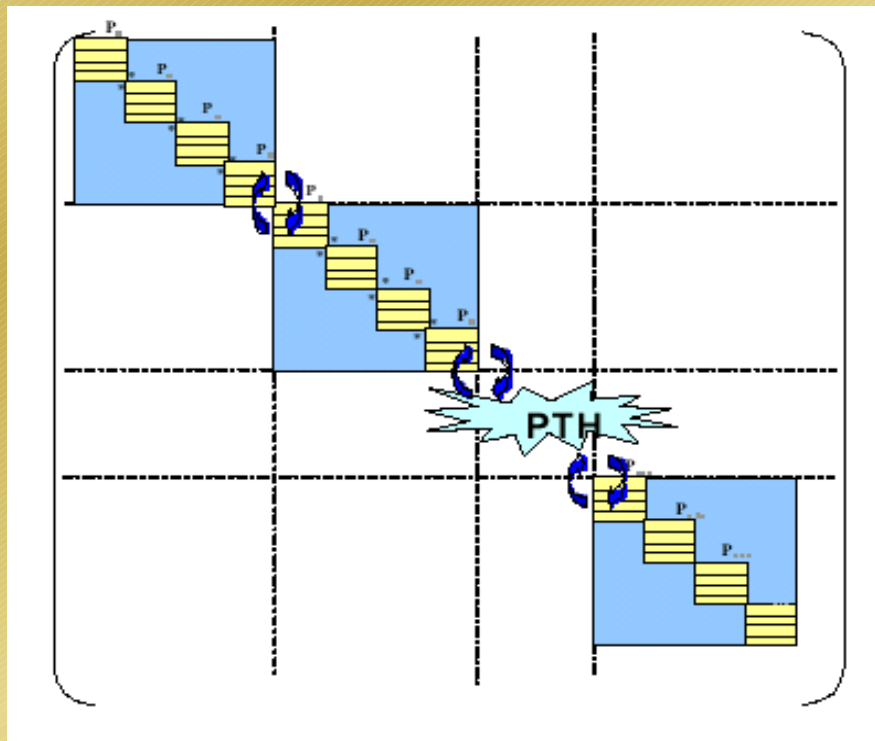
Scaled Speedup of the PDD Algorithm on Paragon. *1024 System of order 1600, periodic & non-periodic*



Scaled Speedup of the Reduced PDD Algorithm on SP2. *1024 System of Order 1600, periodic & non-periodic*

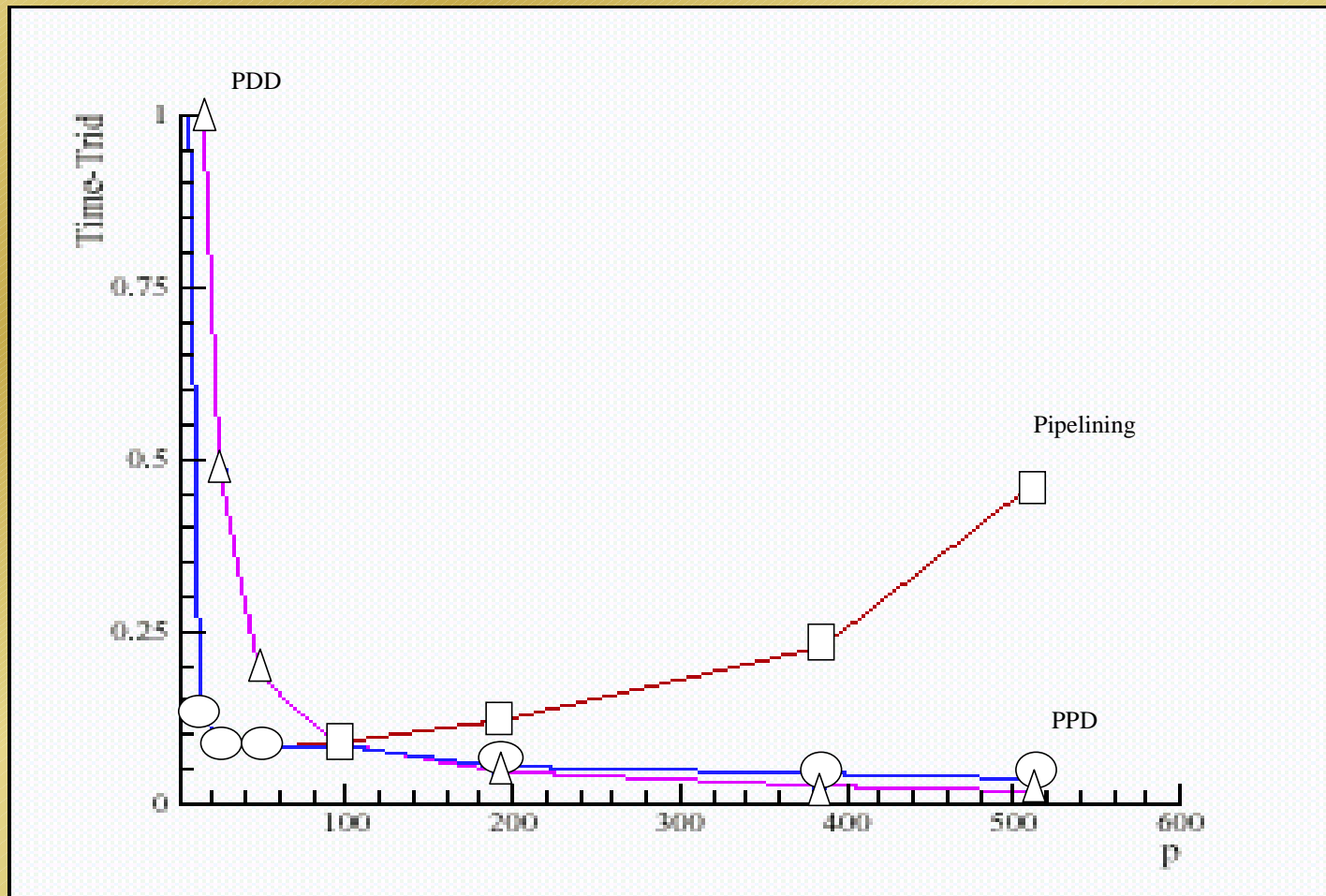
The Parallel Two-Level Hybrid Method

Combine the PDD with the Pipelining

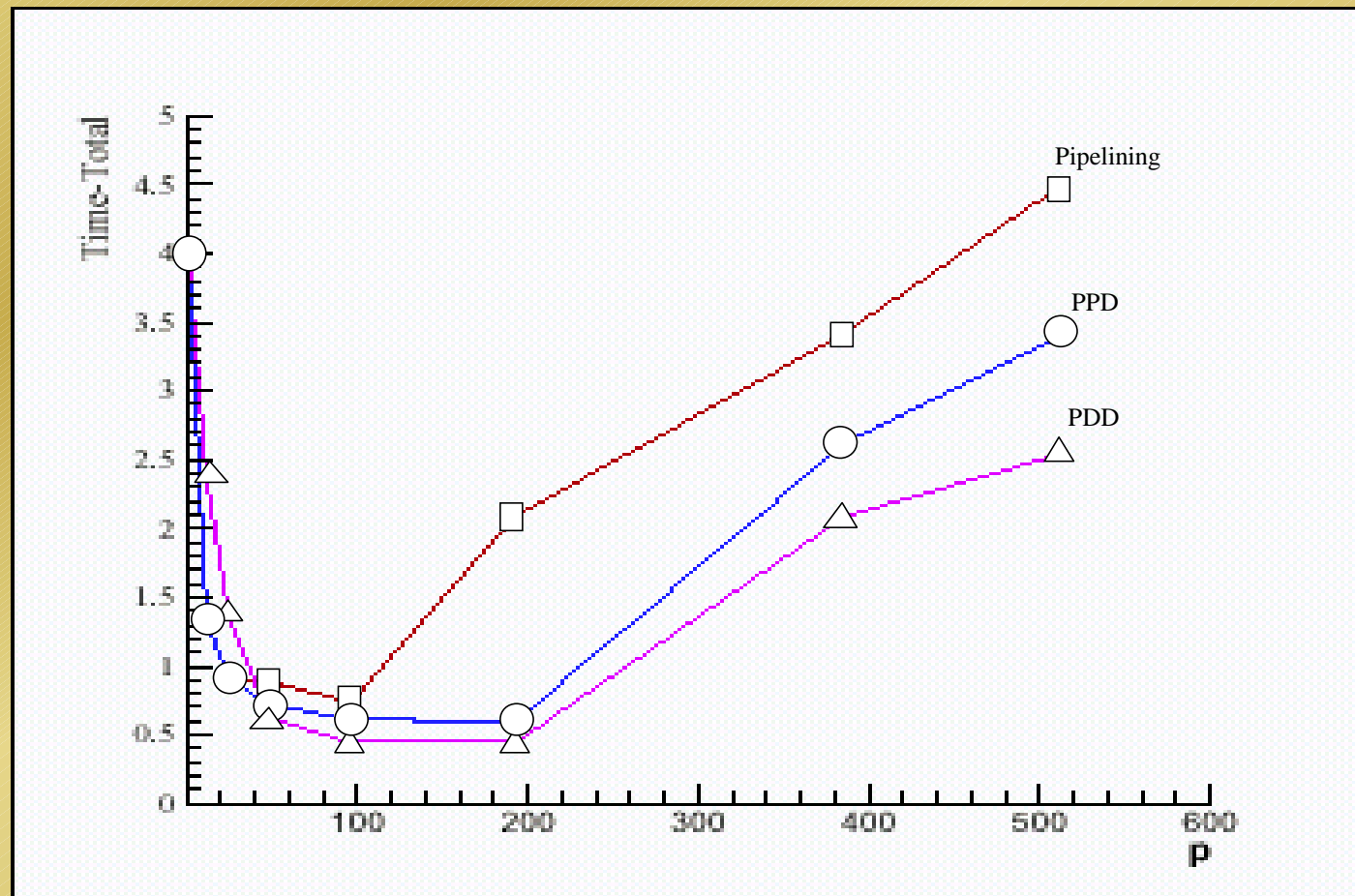


- Use an accurate parallel tridiagonal solver to solve the m super-subsystems concurrently, each with k processors
- Modify PDD algorithm and consider communications only between the m super-subsystems.

PDD, Pipelining, PPD (hybrid) for NLOM Tridiagonal Systems



PDD, Pipelining, PPD for NLOM Poisson Equations



Distributed Computing

High Performance Computing Mobility (HPCM)

- With the rapid advance of communication, the next generation computing will be: **Mobile Computing**
- Current successes of mobile computing are based on safe-languages such as Java, which is slow and cannot apply to legacy codes
- The HPCM middleware makes codes written in traditional languages such as Fortran, C, C++ migratable

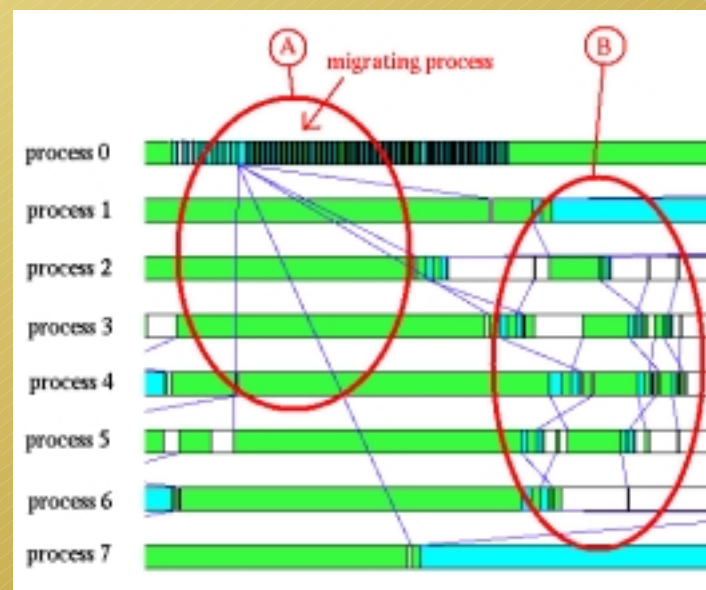
Technical Challenges of HPC mobility

- How to transfer Execution State?
- How to transfer Memory State?
- How to transfer Communication State?
- How to transfer process state efficiently?
- How to transfer process state automatically?
- How to support process migrate from one virtual organization to another in a Grid environment?
- How to support mobility in hybrid Java-native code environment?
- How to design a coordinated [middleware](#)?

Mobility of Legacy Code

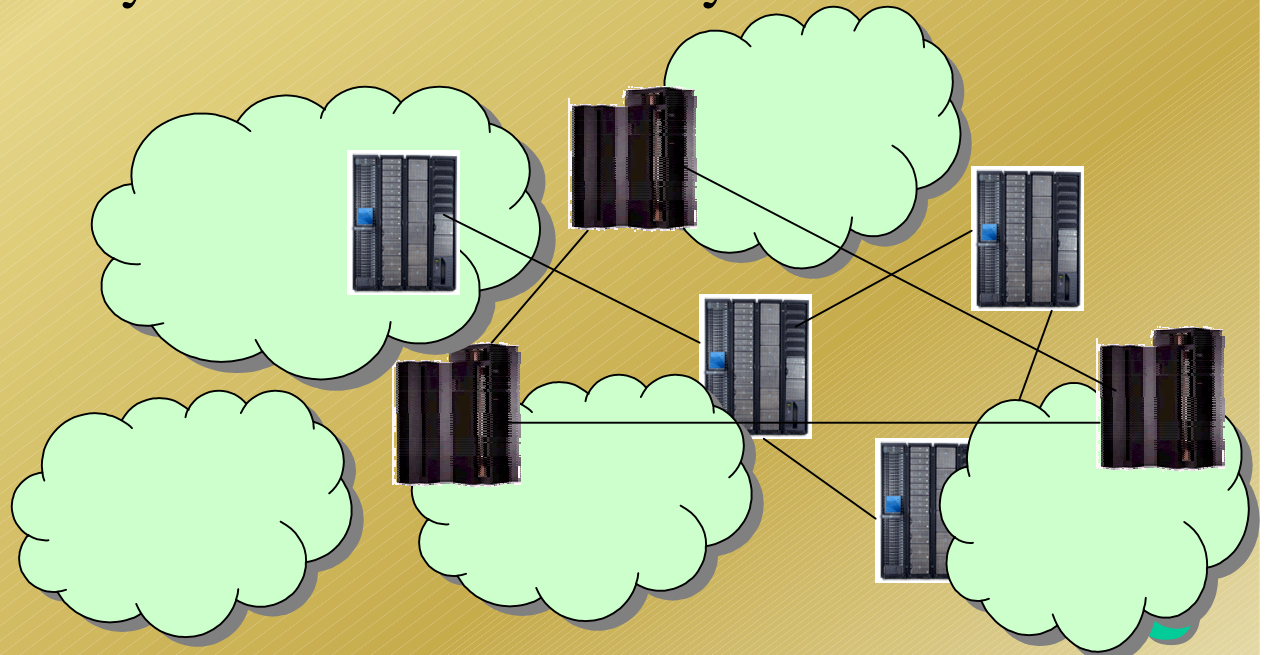
continue

- We have developed novel methodologies and a prototype system, **HPCM**, to migrate codes written in traditional languages such as Fortran, C, C++
 - Two level mobility: migrate native codes under Java virtual machine
 - General methods: migrate between different computing systems and different virtual organizations.
 - Leading technology, strong mobility



Performance and Scheduling

- New challenge in Grid Computing
 - Resources are shared
 - Data are shared resources
- New challenge in high performance computing
 - Memory hierarchy and data access delay

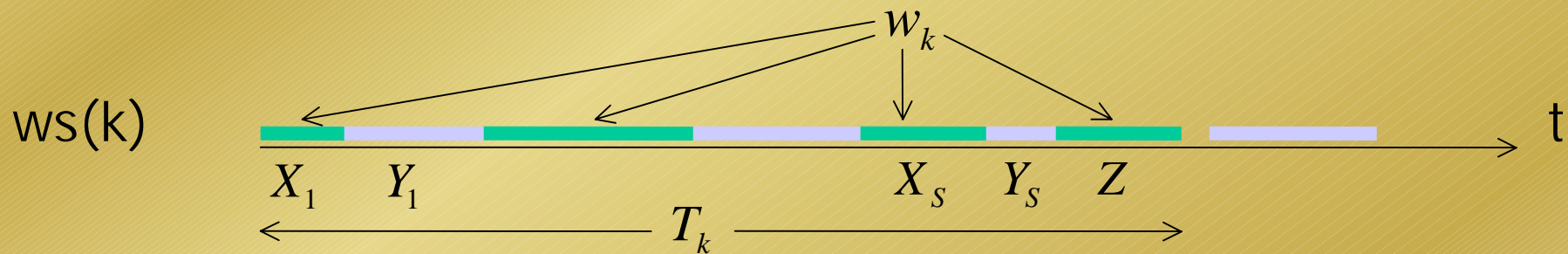


The Grid Harvest Service (GHS) System

- A long-term application-level performance prediction and scheduling system for non-dedicated distributed (Grid) environments
- A new prediction model derived via probability analysis and simulation
- New scheduling heuristics for resource, QoS, and data conscious scheduling
- Runtime dynamic scheduling

Performance Model

- Remote job has low priority
- Local job arriving and service time based on extensive monitoring and observation



$$T_k = X_1 + Y_1 + X_2 + Y_2 + \dots + X_s + Y_s + Z$$

$$T_k = w + Y_1 + Y_2 + \dots + Y_s$$

Prediction Formula

- Parallel task completion time

$$\Pr(T \leq t) = \begin{cases} \prod_{k=1}^m [e^{-\lambda_k w_k} + (1 - e^{-\lambda_k w_k}) \Pr(U(S_k) \leq t - w_k \mid S_k > 0)], & \text{if } t \geq w_{\max} \\ 0, & \text{otherwise} \end{cases}$$

- Homogeneous parallel task completion time

$$\Pr(T \leq t) = \begin{cases} [e^{-\lambda w} + (1 - e^{-\lambda w}) \Pr(U(S) \leq \tau \mid S > 0)]^m, & \text{if } \tau > 0 \\ 0, & \text{otherwise} \end{cases}$$

where, $\tau = t - w$

- Mean time balancing partition

$$w_k = \frac{W}{\sum_{k=1}^m (1 - \rho_k) \tau_k} (1 - \rho_k) \tau_k$$

Scheduling Algorithms

Scheduling with a Given Number of Sub-tasks

List a set of lightly loaded machines $M = \{m_1, m_2, \dots, m_q\}$;

List all possible sets of machines, such as $|S_i| = p$

For each machine set S_k ($1 \leq k \leq z$),

Use **mean time balancing partition** to partition the task

Use the formula to calculate the mean and coefficient of variation

If $E(T_{S_{p'}})(1 + Coe.(T_{S_{p'}})) > E(T_{S_k})(1 + Coe.(T_{S_k}))$, then $p' = k$;

End For

Assign parallel task to the machine set $S_{p'}$;

Optimal Scheduling Algorithm

List a set of lightly loaded machines $M = \{m_1, m_2, \dots, m_q\}$;

While $p < q$ **do**

Scheduling with p Sub-tasks

If $E(T_{S_{k'}, p'}) (1 + Coe.(T_{S_{k'}, p'})) > E(T_{S_k, p}) (1 + Coe.(T_{S_k, p}))$, then

$p' = p$;

End If

End while

Assign parallel task to the machine set $S_{k'}^{p'}$.

Heuristic Scheduling Algorithm

- List a set of lightly loaded machines $M = \{m_1, m_2, \dots, m_q\}$;
- Sort the machines in a decreasing order with $(1 - \rho_k)\tau_k$;
- Use the **task ratio** to find the upper limit q ;
- Use bi-section search to find the p such as

$$E(T_{S_k^p})(1 + Coe.(T_{S_k^p}))$$

is minimum

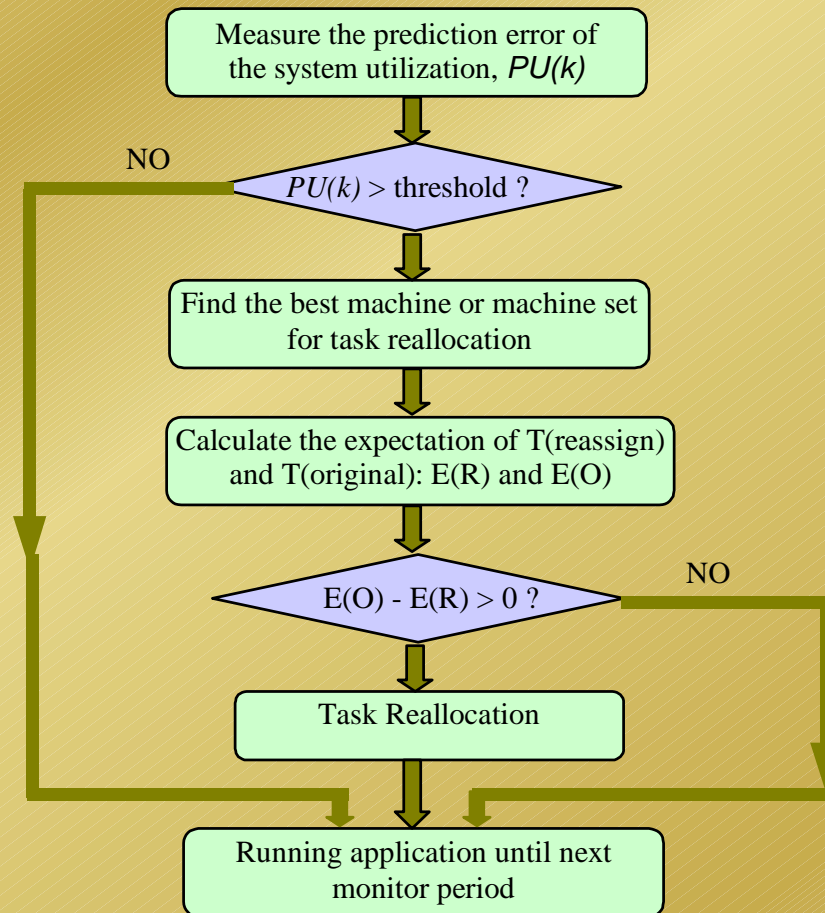
QoS Guided Min-Min Heuristics

```
for all tasks  $t_i$  in meta-task  $M_v$  (in an arbitrary order)
  for all hosts  $m_j$  (in a fixed arbitrary order)
     $CT_{ij} = ET_{ij} + d_j$ 
  do until all tasks with high QoS request in  $M_v$  are mapped
    for each task with high QoS in  $M_v$ , find a host in the QoS qualified host
    set-          that obtains the earliest completion time
    find the task  $t_k$  with the minimum earliest completion time
    assign task  $t_k$  to the host  $m_l$  that gives it the earliest completion time
    delete task  $t_k$  from  $M_v$ 
    update  $d_l$ 
    update  $CT_{il}$  for all  $i$ 
  end do
do until all tasks with low QoS request in  $M_v$  are mapped
  for each task in  $M_v$  find the earliest completion time and the
  corresponding host
  find the task  $t_k$  with the minimum earliest completion time
  assign task  $t_k$  to the host  $m_l$  that gives it the earliest completion
  time
  delete task  $t_k$  from  $M_v$ 
  update  $d_l$ 
  update  $CT_{il}$  for all  $i$ 
end do
```

Data-Conscious Scheduling Heuristics

```
For a group of tasks  $T_i$ 
  Get  $MCT_0$  based on data replica placement.
For each task in the metatask but not run yet
  If subtask  $t_{ij}$  needs data  $d_n$ 
    For all  $S_m$  in Sites
      Compute  $MCT_m$  based on  $DDT$  from host 1 to  $m$ .
      Get minimum  $MCT_{min}$ 
    Endfor
    If  $MCT_{min} < MCT_0$ 
      Copy data  $j$  from host 1 to host  $r$ 
       $MCT_0 = MCT_{min}$ 
    Else
      Do not copy
  Endfor
```


Rescheduling Algorithm



GHS Design: System Architecture

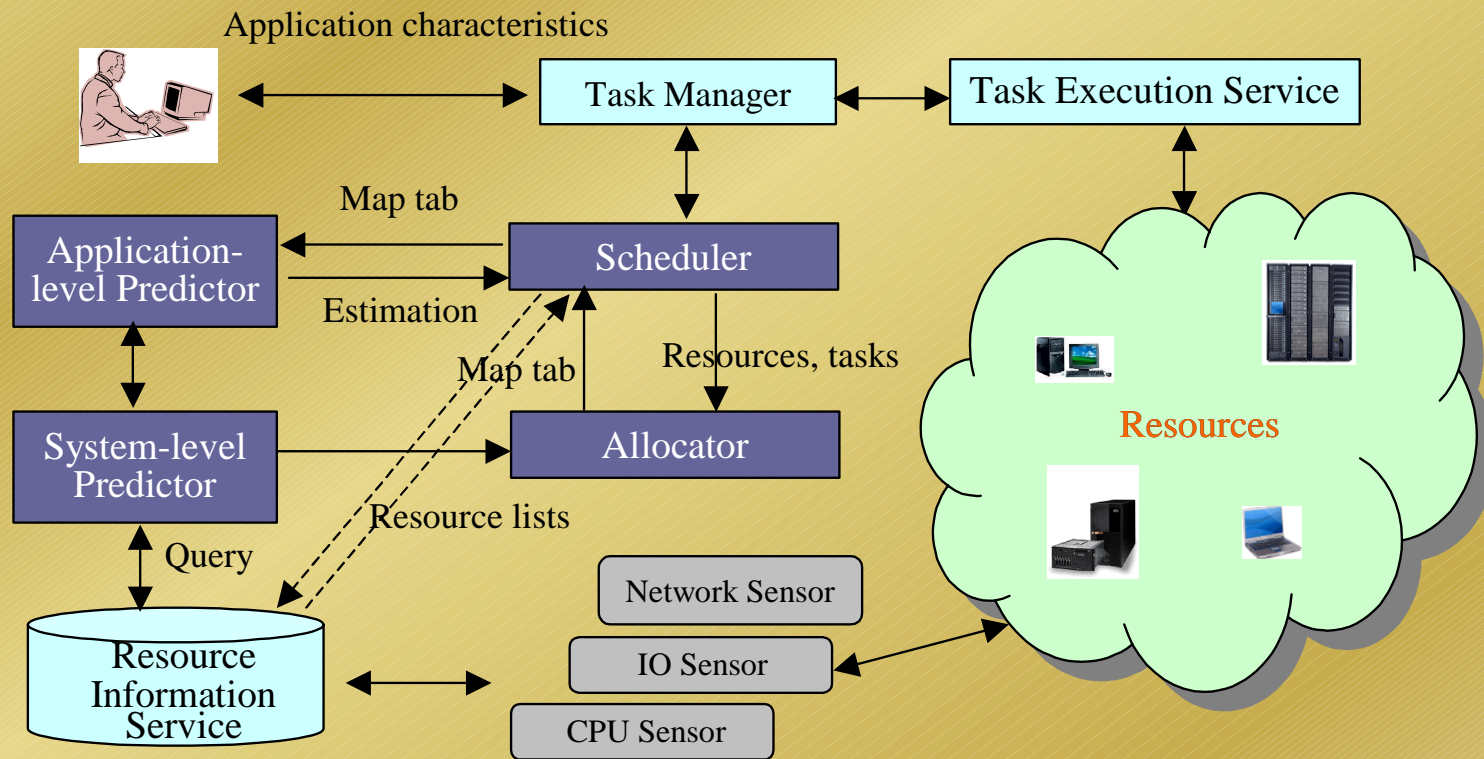
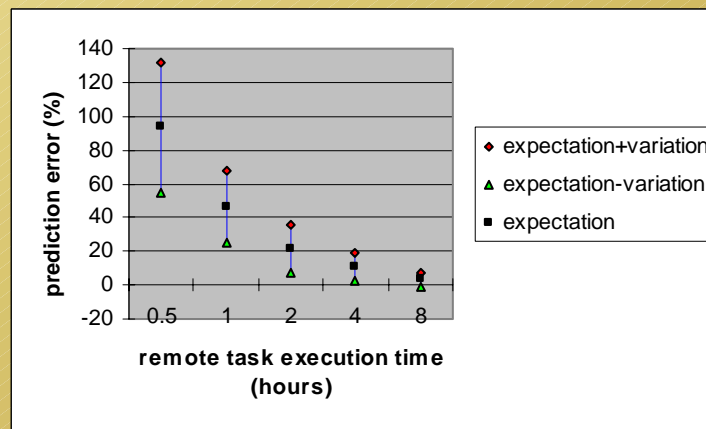


Figure. 1. A framework of GHS task scheduling system

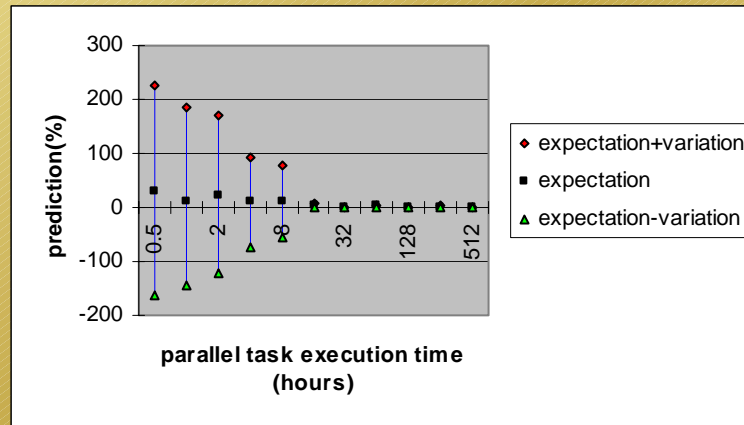
Experimental Testing

Application-level Prediction

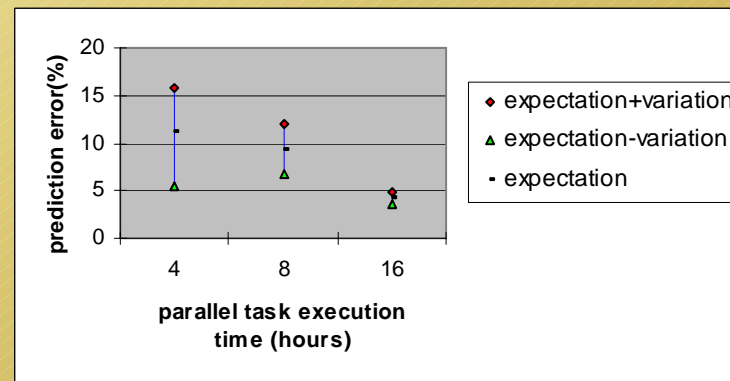
$$\left| \frac{\text{Prediction}_{\text{period}} - \text{Measurement}}{\text{Measurement}} \right|$$



Remote task completion time on single machine

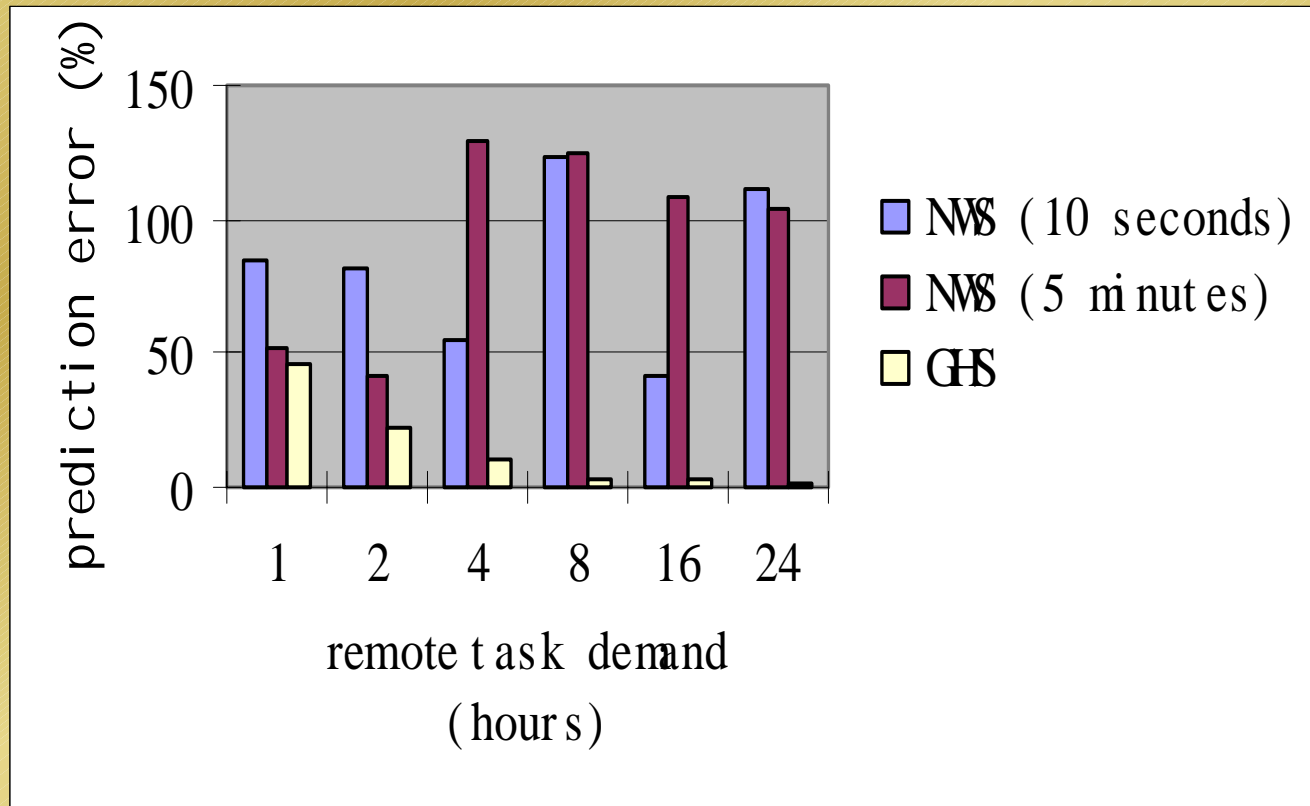


Prediction of parallel task completion time



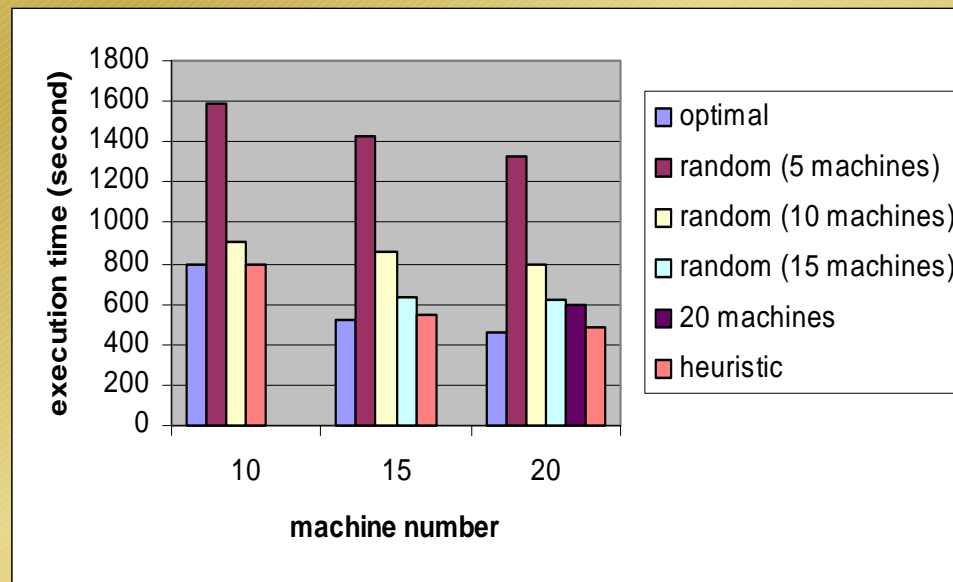
Prediction of a multi-processor with local scheduler

Comparison with NWS



Mean of the prediction error of NWS and GHS

Performance Gain with Scheduling



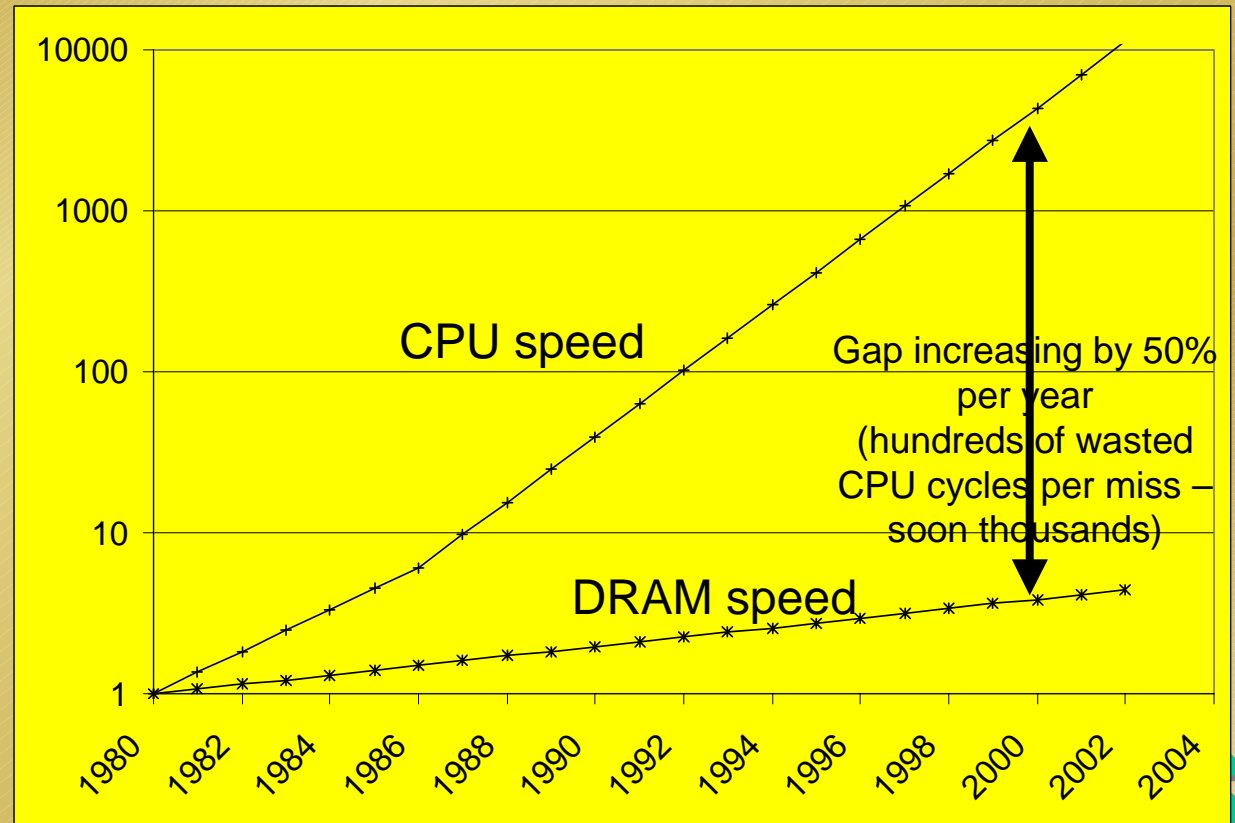
Execution time with different scheduling strategies

Data Access Optimization

- Data access is the bottleneck of high performance computing

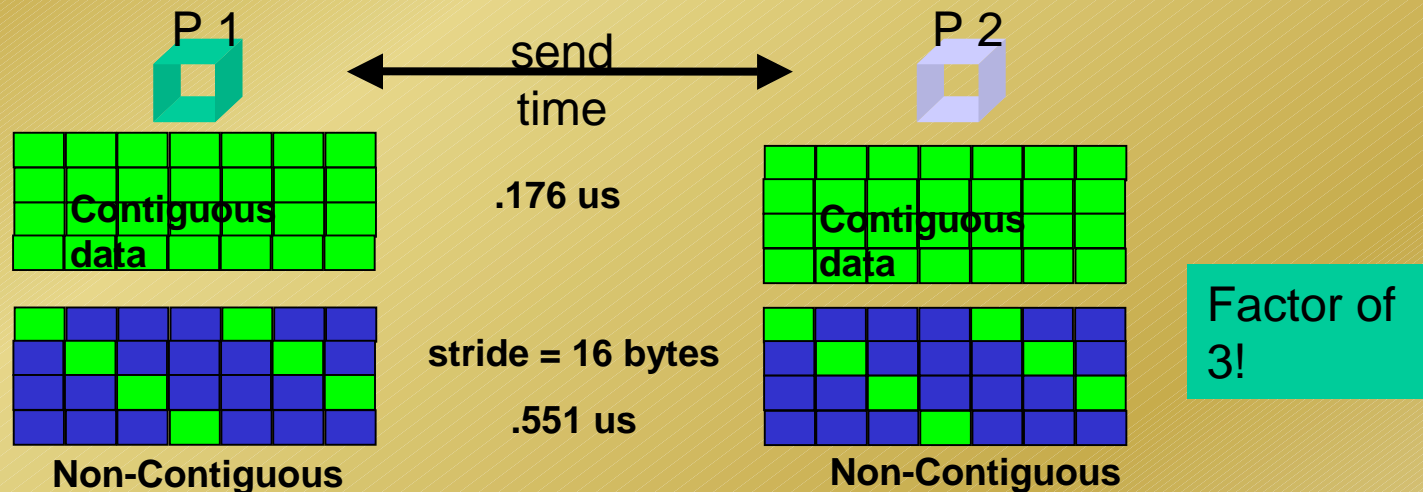
Relative Speeds of CPU verse DRAM

Much worse for I/O
bottleneck
(millions of wasted
CPU cycles per
miss)



Solution Approach

- Memory access pattern significant impact on comm.
- Impact differently on different machines
- The memory LogP model
- Optimization methods
- Application on MPI applications



Memory communication cost can be significant for real applications.

The memory logP model

Memory hierarchy

CPU

Sender user space

Application buffer

Network Buffer

System space
(Interconnect Cloud)

Network Buffer

Receiver user space

Application buffer

CPU

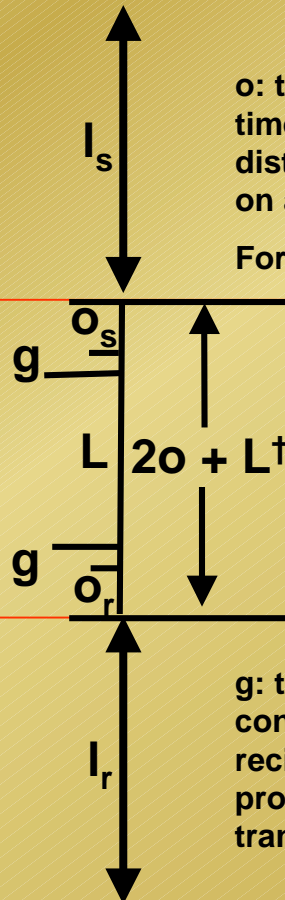
X. Sun, F. ...

Memory hierarchy

l : the additional latency in transfer non-contiguous data of data size (s) and distribution (d) for a given implementation of data transfer on a given system, $l=f(s,d)$

o : the data transfer overhead, defined as the length of time to transfer contiguous data of data size (s) and distribution (d) for a given implementation of data transfer on a given system, $l=f(s,d)$

For ideal case, o remains constant with size.



$$\text{Total Comm Cost} = 2(o+l) + L$$

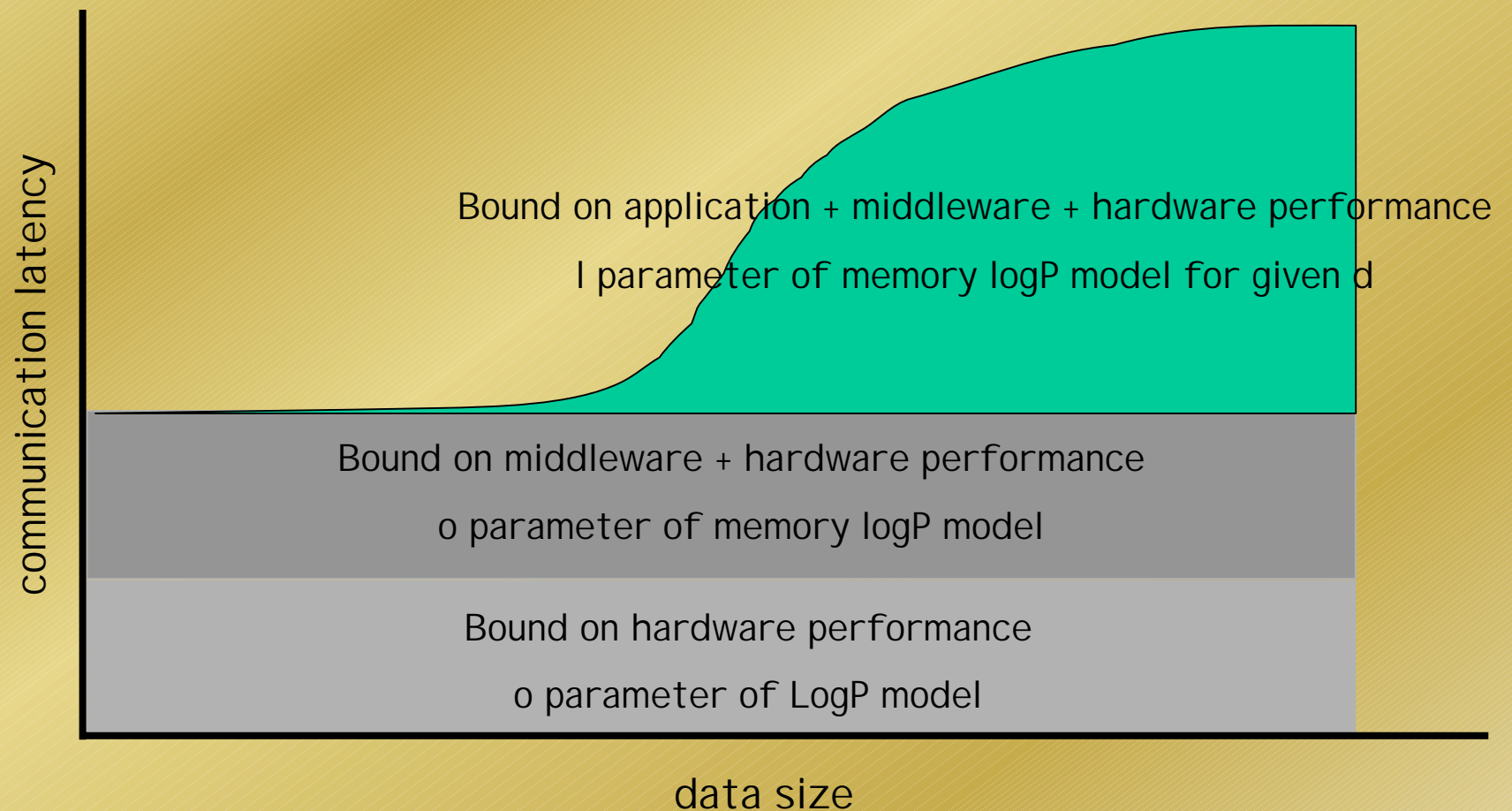
g : the gap, defined as the minimum time interval between consecutive message receptions at a processor. The reciprocal of g corresponds to the available per-processor bandwidth for a given implementation of data transfer on a given system.

P : the number of processor/memory modules

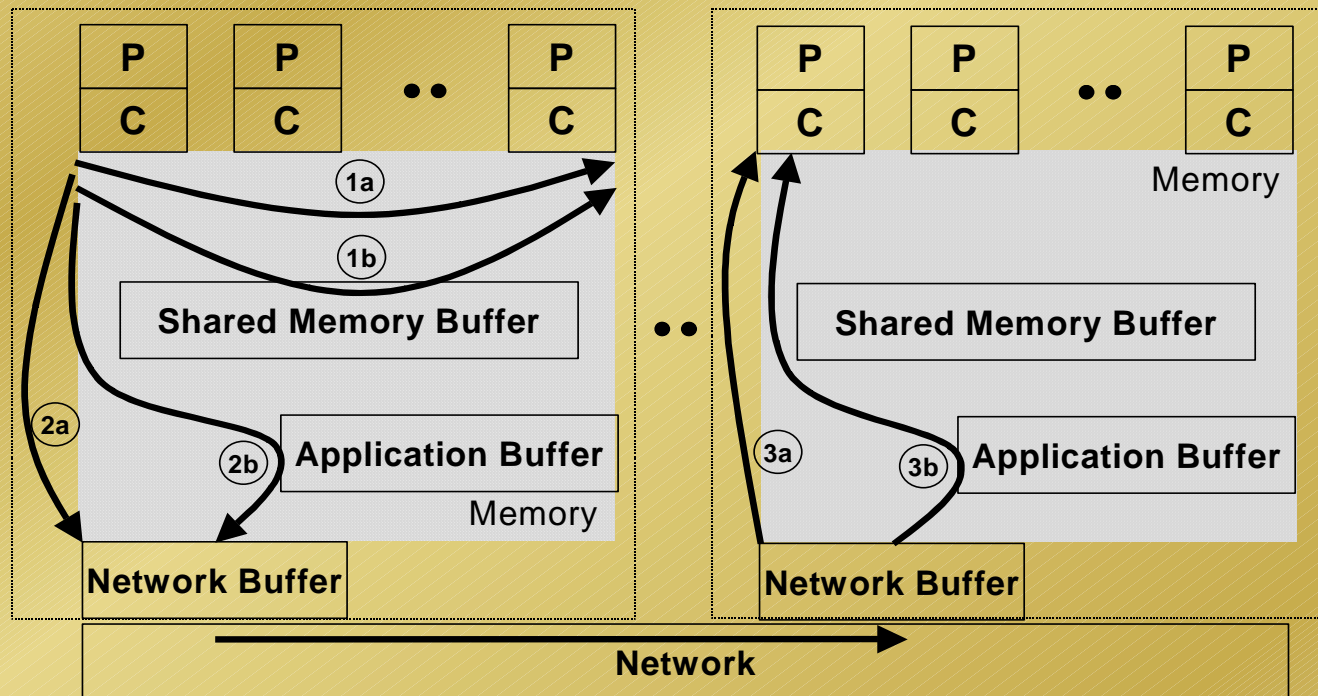
† Assume $g = o$, common simplification of LogP model



Performance bounds using memory-logP

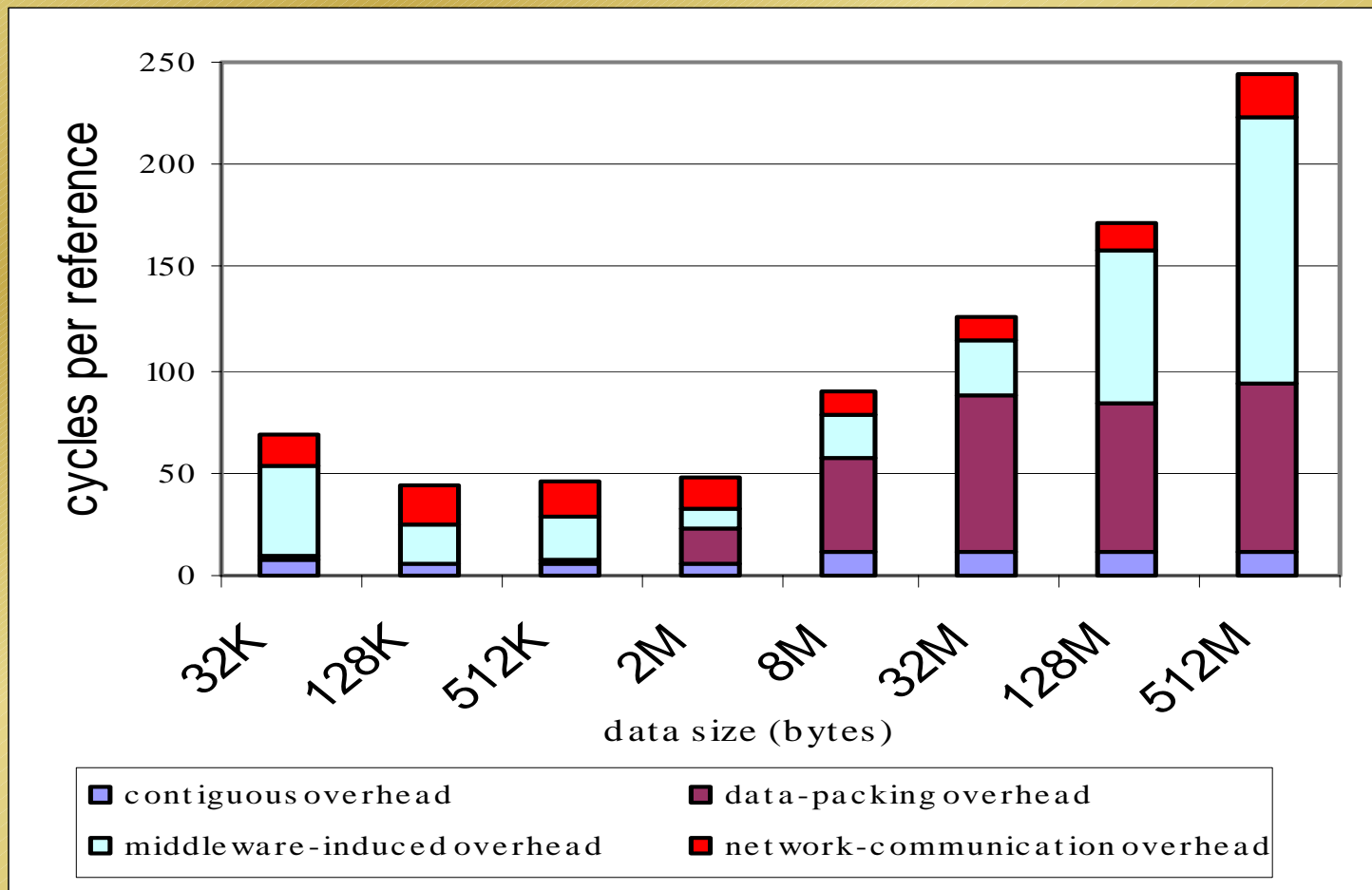


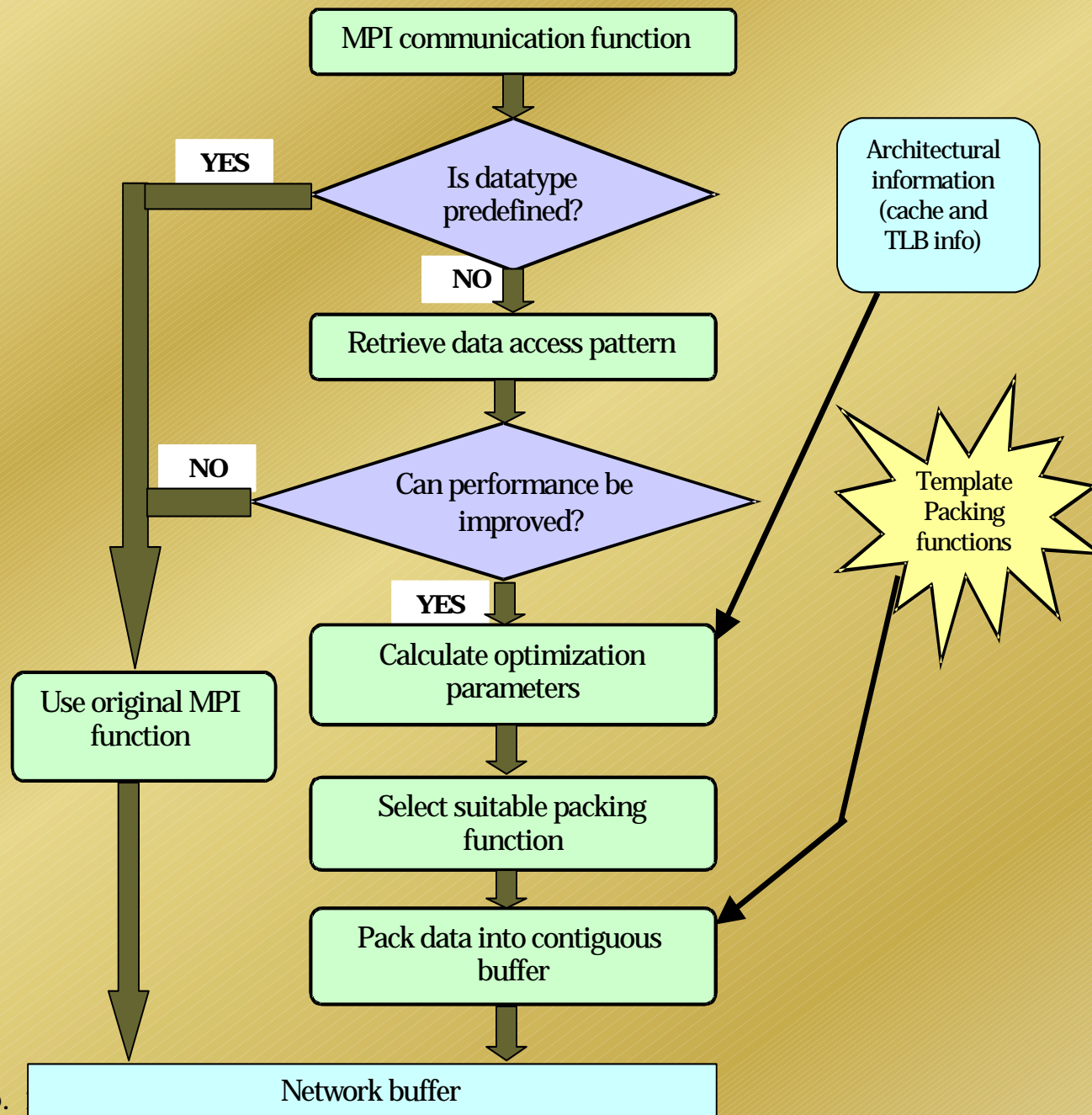
Case Study: Matrix Transpose on a SGI O2000 using MPI middleware for communication



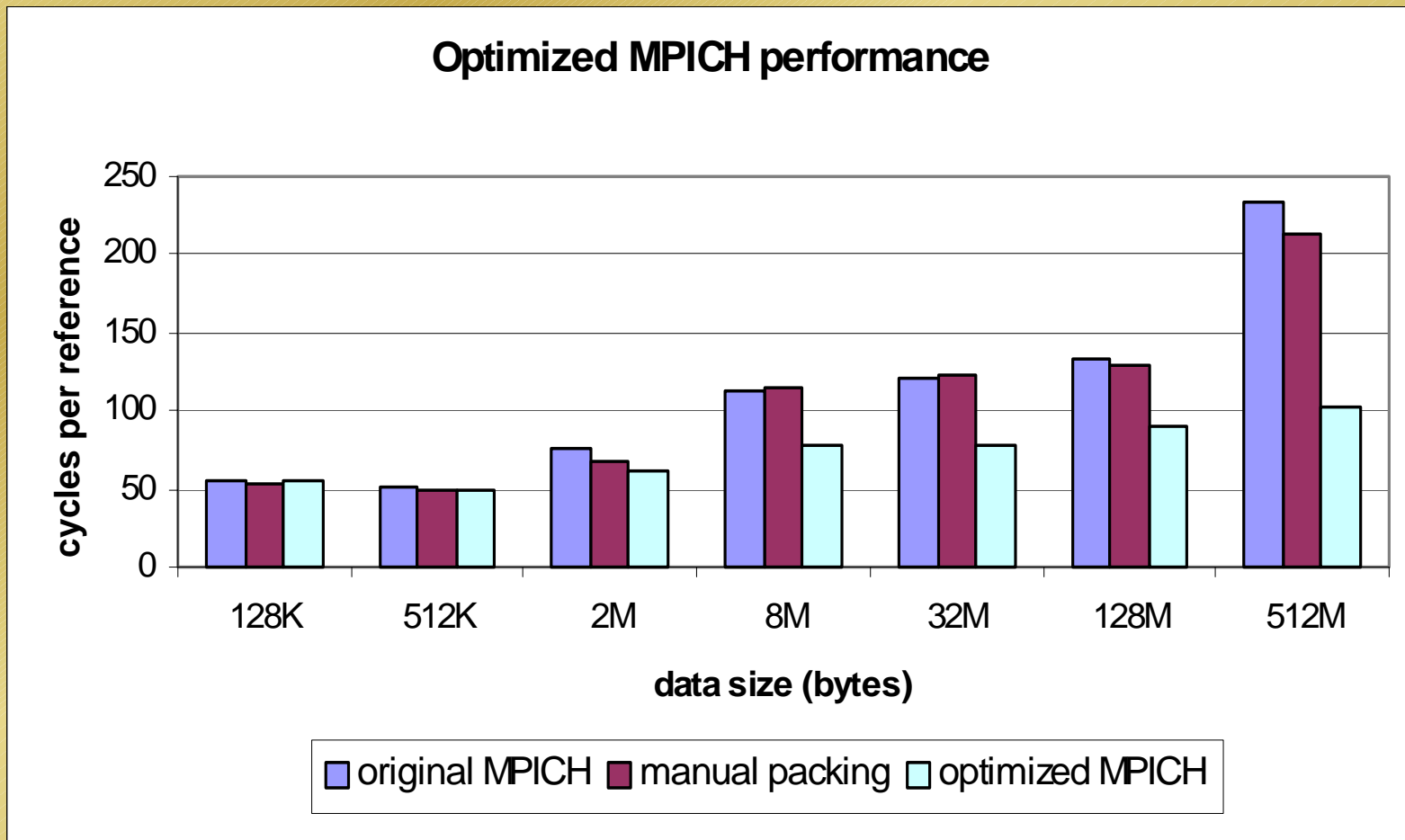
Comm. critical paths: Non-Cont. data need
go through the application buffer first

Quantifying Communication Cost for Matrix Transpose

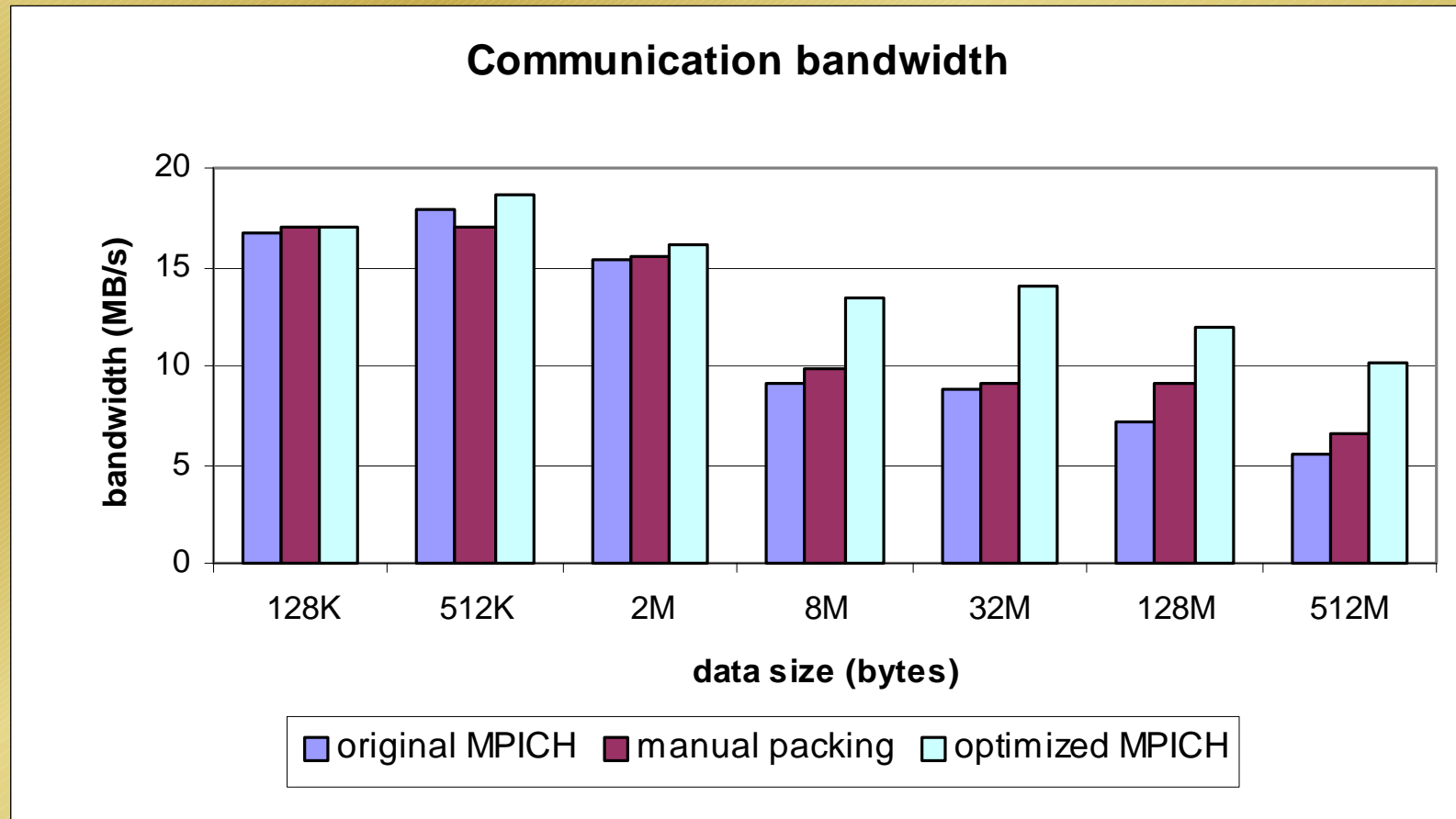




Improved performance of MPI derived datatypes



Bandwidth Improvement



The Cosmology Application

- Worked on a cosmology application called ENZO
 - SAMR (Structured Adaptive Mesh Refinement) application in astrophysics and cosmology
 - Developed by G. Bryan and M. Norman
- Three major tasks:
 - Conducted a detailed performance analysis
 - Developed a novel dynamic load balancing (DLB) scheme for ENZO on parallel systems
 - Extended the code to the Computational Grid

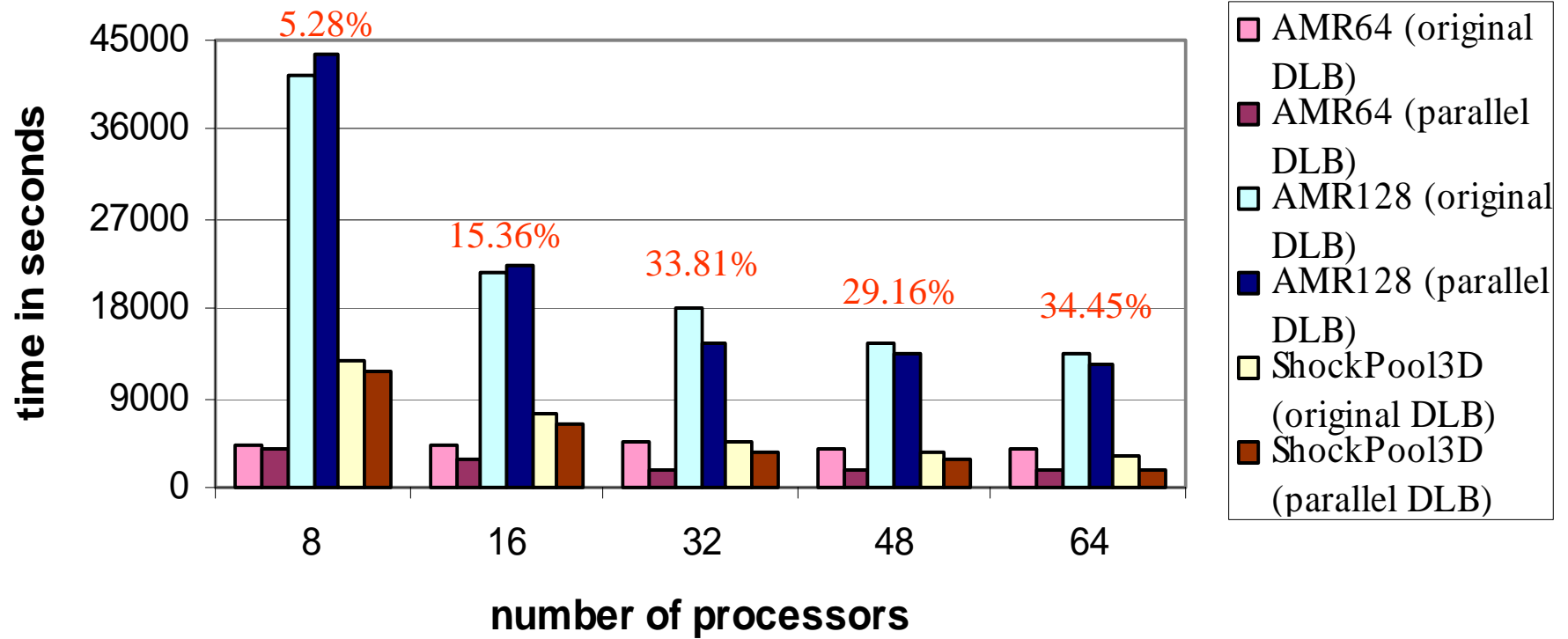
Task 1: Performance Analysis

- Overall characteristics
 - Manual instrumentation with FPMPI tool
- I/O performance
 - Use Pablo toolkit
- Adaptive characteristics
 - Coarse granularity
 - High magnitude of imbalance
 - Different patterns of imbalance
 - High frequency of refinements

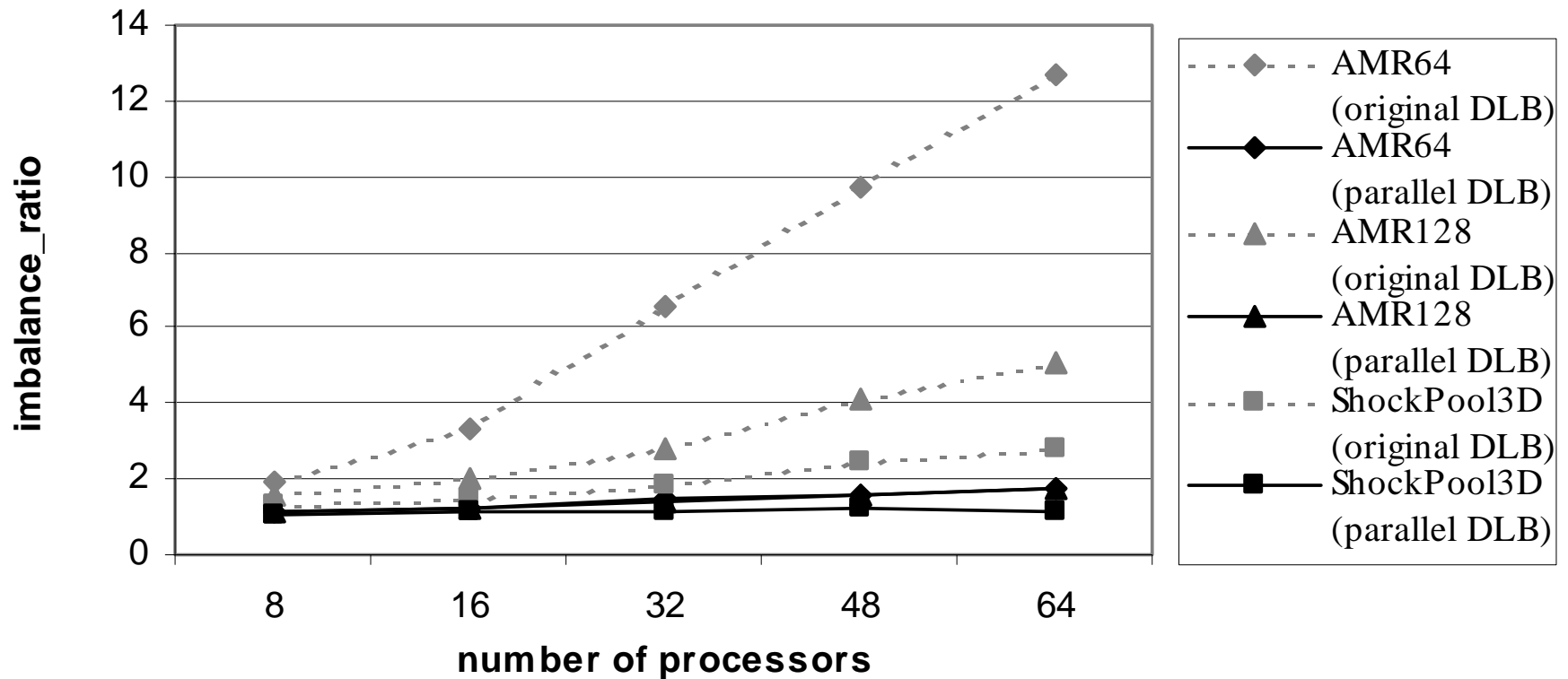
Task 2: Parallel DLB

- Design a DLB scheme for SAMR applications considering their adaptive characteristics
- Moving-grid phase
 - Directly move excess grids from overloaded proc. to underloaded proc.
 - Minimize grid movements by global information
 - Address the high frequency and different patterns
- Splitting-grid phase
 - Split the largest grid on the overloaded proc.
 - Address the coarse granularity

Total Execution Time



Imbalance Ratio



$$imblanace_ratio = \frac{MaxLoad}{AvgLoad}$$

Task 3: Extension to Grid

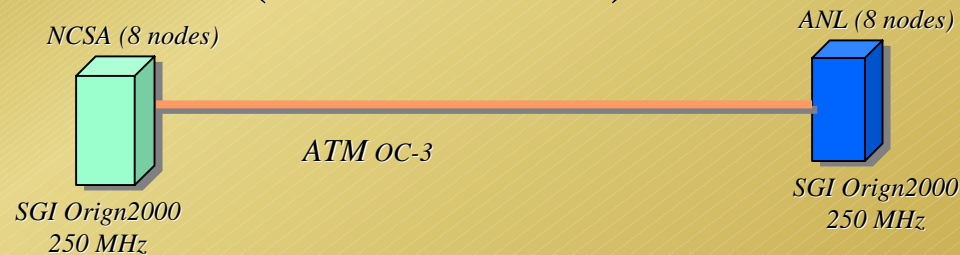
- Extend the improved version to the Grid
 - Alliance's Grid
- One major issue: a new DLB scheme for distributed environment
 - Heterogeneity of processors
 - Heterogeneity of networks
 - Dynamic features of networks
 - Adaptive features of applications

Experiments

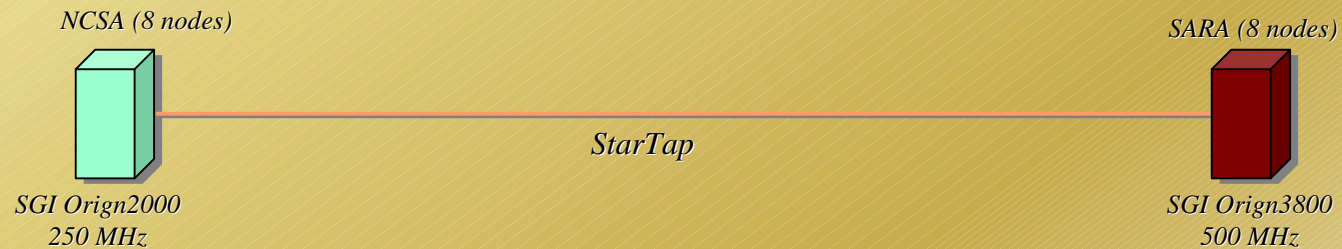
- DS1: LAN-connected



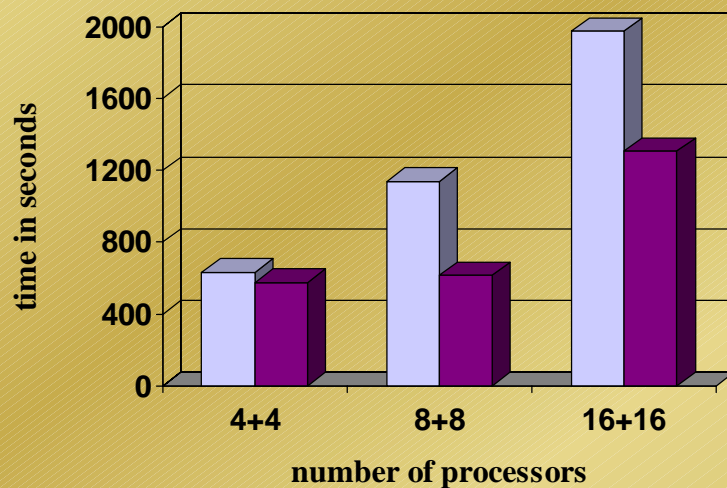
- DS2: WAN-connected (two locations)



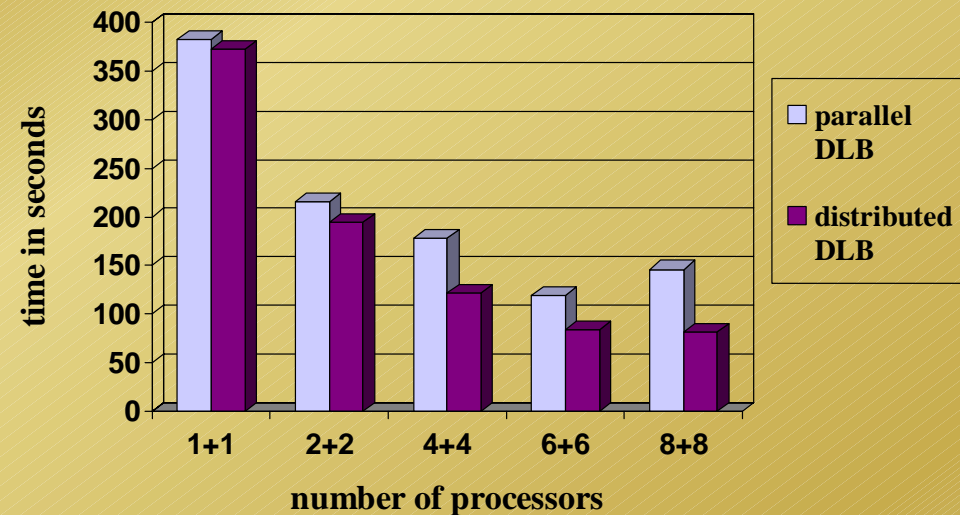
- DS3: WAN-connected (two countries)



Execution Time for AMR64 on DS1

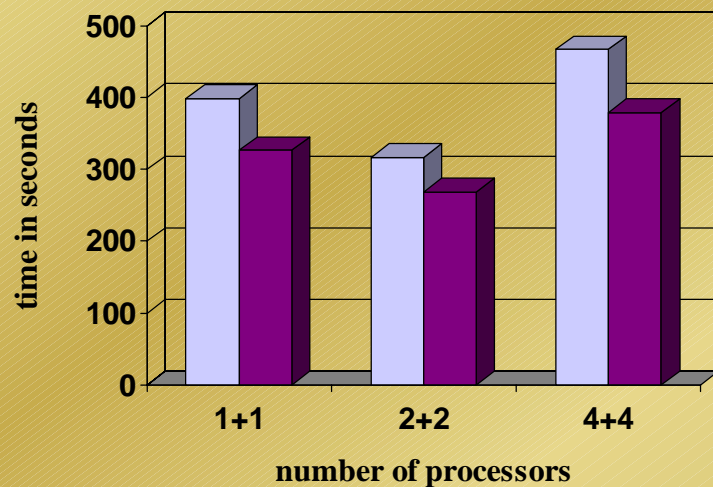


Execution Time for ShockPool3D on DS2

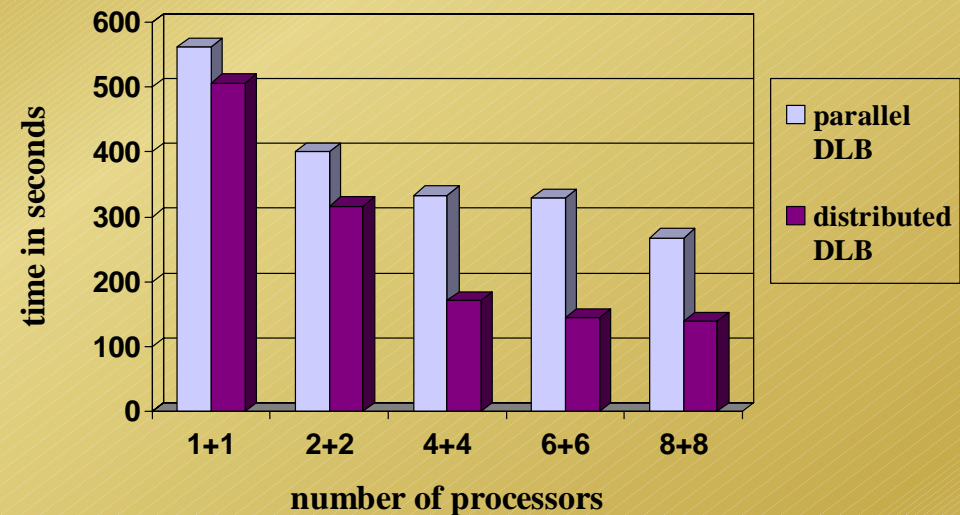


👉 The relative improvement ranges from 2.6% to 45.9%

Execution Time for AMR64 on DS3



Execution Time for ShockPool3D on DS3



👉 The relative improvement ranges from 10.0% to 56.1%

Conclusion

- **Memory Model and Algorithms:** Improve high performance computing
- **Mobility plus Dynamic Scheduling:** improve performance, reliability, availability, QoS, and trustiness of distributed computing
- **Application-level** Performance Optimization
- **Fermi Applications:** Performance analysis and enhancement, implementation on Grid, such as the TeraGrid